

# Computer Vision

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله  
للعلوم والتقنية  
King Abdullah University of  
Science and Technology



KAUST Academy  
King Abdullah University of Science and Technology

January 28, 2026

1. Introduction.
2. Applications of Computer Vision.
3. Image Representation.
4. Convolutional Neural Network (CNN) and its Components.
5. CNN-based Architectures (AlexNet, VGG, InceptionNet, ResNet, EfficientNet, and MobileNet).

- ▶ Understand the basic concepts of Computer Vision and its real-world applications.
- ▶ Describe how images are represented and processed in a computer.
- ▶ Explain the fundamental building blocks of Convolutional Neural Networks (CNNs).
- ▶ Differentiate between popular CNN architectures (AlexNet, VGG, InceptionNet, ResNet, EfficientNet, and MobileNet) and their key innovations.

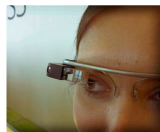
Building artificial systems that process, perceive, and reason about visual data

# Computer Vision is Everywhere

→ lab samples  
→ CCTV  
→ inspection  
→ Defects



Left to right:  
[image](#) is [CC BY-SA](#) public domain  
[image](#) is [CC BY-SA](#) public domain  
[image](#) by [NASA](#) is licensed under [CC BY-SA](#)  
[image](#) is [CC BY-SA](#) public domain



Bottom row, left to right:  
[image](#) is [CC BY-SA](#) public domain  
[image](#) by [Google Glass](#) is licensed under [CC BY-SA](#) changes made  
[image](#) is public domain  
[image](#) is licensed under [CC BY-SA](#) changes made

# Some Applications

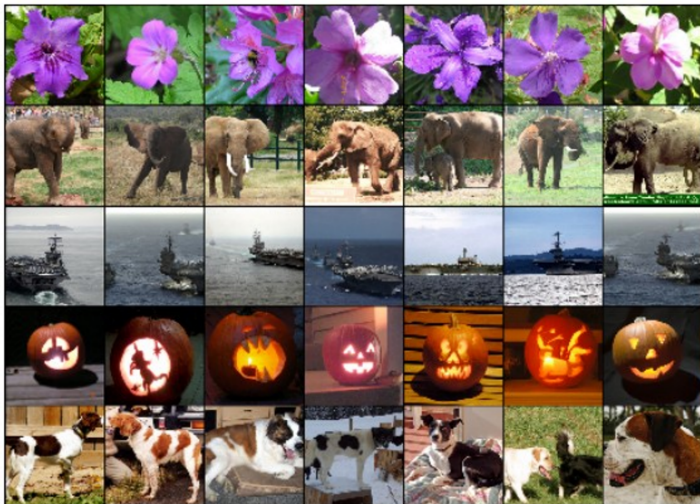
defect/



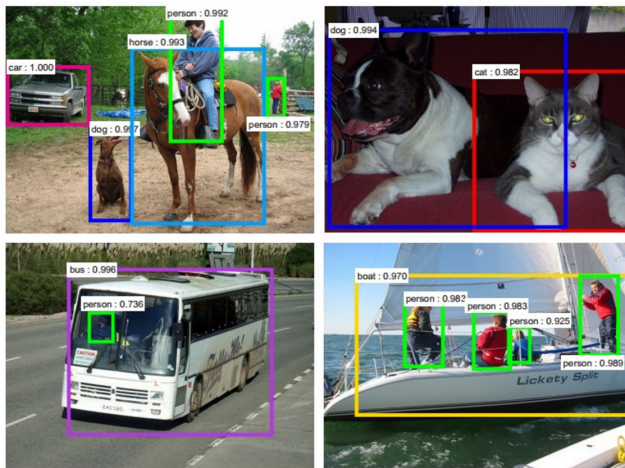
→ AI → Image Classification



## Image Retrieval

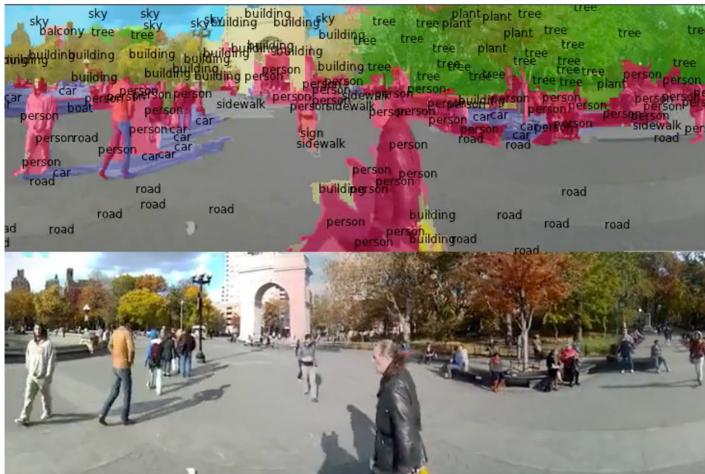


## Object Detection



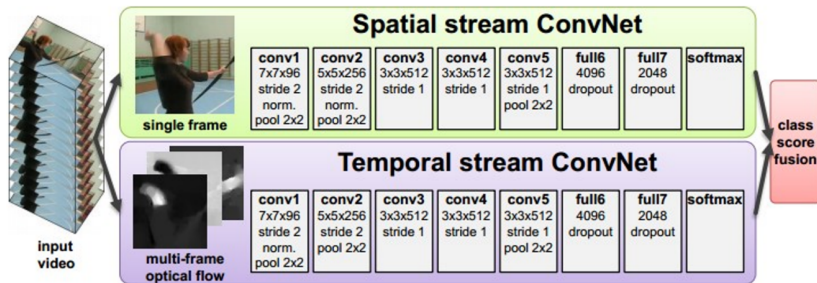
Ren, He, Girshick, and Sun, 2015

## Image Segmentation



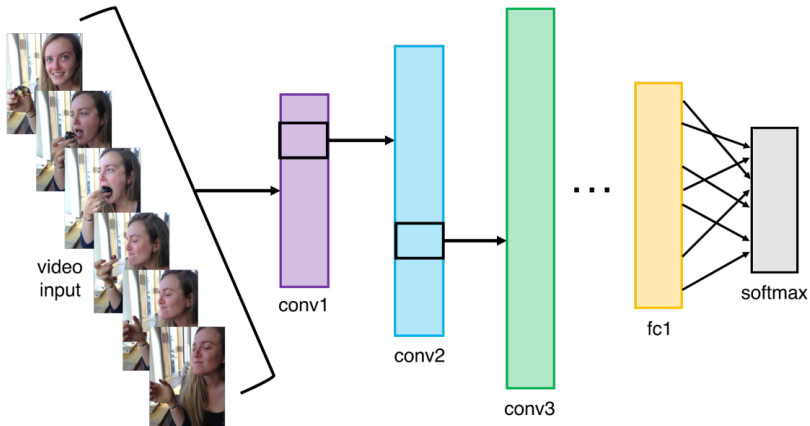
Fabret et al, 2012

## Video Classification



Simonyan et al, 2014

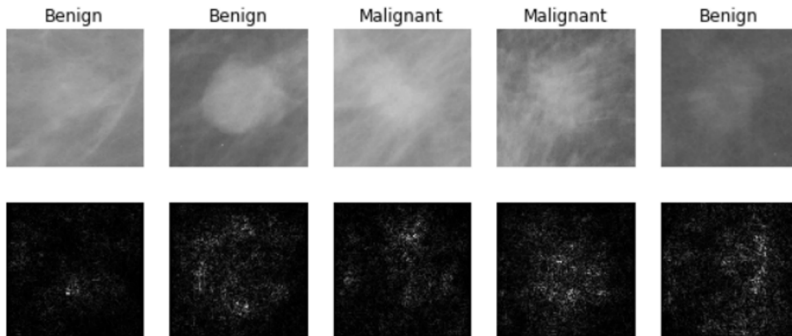
## Activity Recognition



Pose Recognition (Toshev and Szegedy, 2014)



## Medical Imaging





## Image Generation



“Teddy bears working on new AI research underwater with 1990s technology”

DALL-E 2



Style Transfer

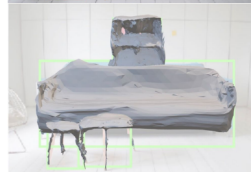
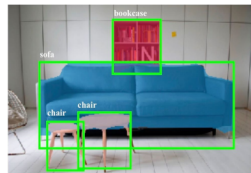
## 3D Vision



Choy et al., 3D-R2N2: Recurrent Reconstruction Neural Network (2016)



Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)



Gkioxari et al., "Mesh R-CNN", ICCV 2019

# How to represent an image?



- ▶ Images are represented as Matrices with elements in  $[0, 255]$
- ▶ Grayscale images have one channel.



187	183	174	168	150	162	129	77	172	161	155	166
155	182	163	74	75	62	93	17	110	210	180	154
180	180	90	14	54	6	10	33	48	106	159	181
206	109	5	134	131	111	120	204	166	15	56	180
194	68	137	281	237	239	239	238	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	181	183	158	227	178	143	182	106	36	190
205	174	195	282	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	168	227	210	127	102	36	101	285	224
190	214	173	66	103	143	56	86	2	108	249	215
187	196	236	75	1	81	47	0	6	217	285	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	95	218

157	183	174	168	150	162	129	161	172	161	155	166
155	182	163	74	75	62	93	17	110	210	180	154
180	180	90	14	54	6	10	33	48	106	159	181
206	109	5	134	131	111	120	204	166	15	56	180
194	68	137	281	237	239	239	238	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	181	183	158	227	178	143	182	106	36	190
205	174	195	282	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	168	227	210	127	102	36	101	285	224
190	214	173	66	103	143	56	86	2	108	249	215
187	196	236	75	1	81	47	0	6	217	285	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	95	218

<sup>0</sup><https://www.v7labs.com/blog/image-recognition-guide>

# How to represent an image?

- ▶ RGB images have 3 channels.



# How can we build vision models?

Input size  $\rightarrow 64 \times 64 \times 3$

- ▶ So far, we've worked with tabular data, where each sample consists of structured features. We processed this data using Fully Connected Neural Networks.

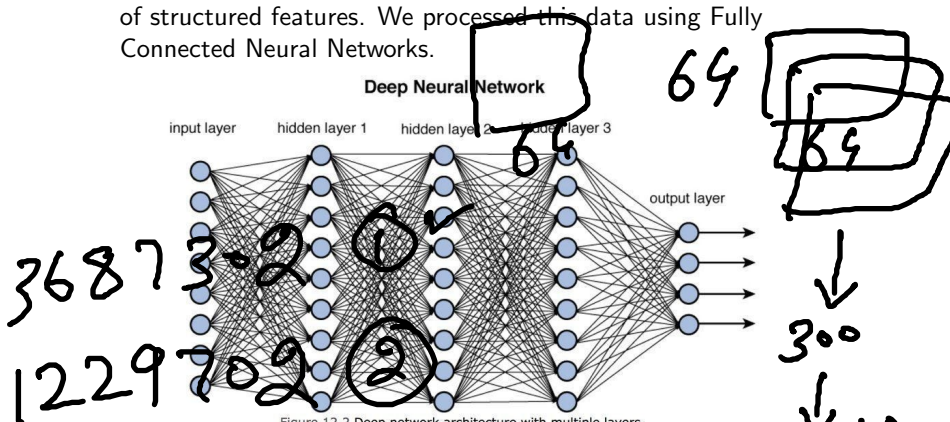
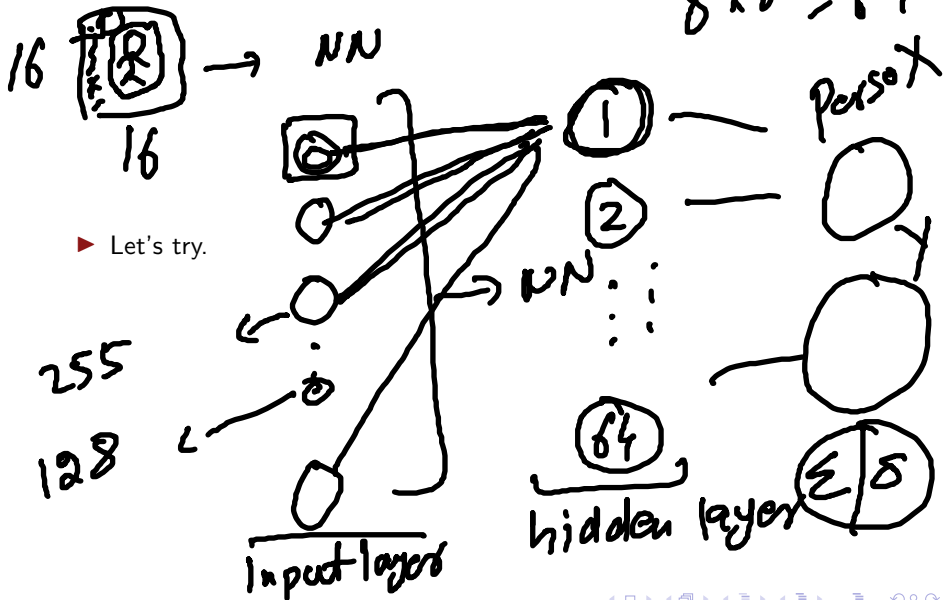


Figure 12.2 Deep network architecture with multiple layers.

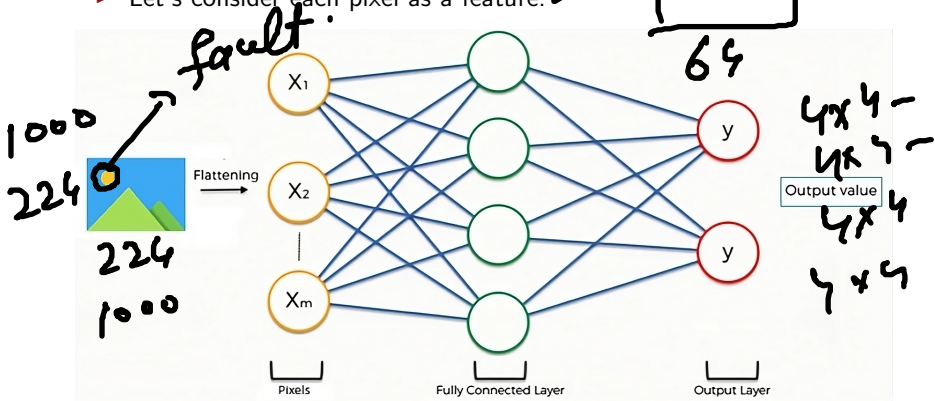
- ▶ But can we apply the same approach to images?

# How can we build vision models?



# How can we build vision models?

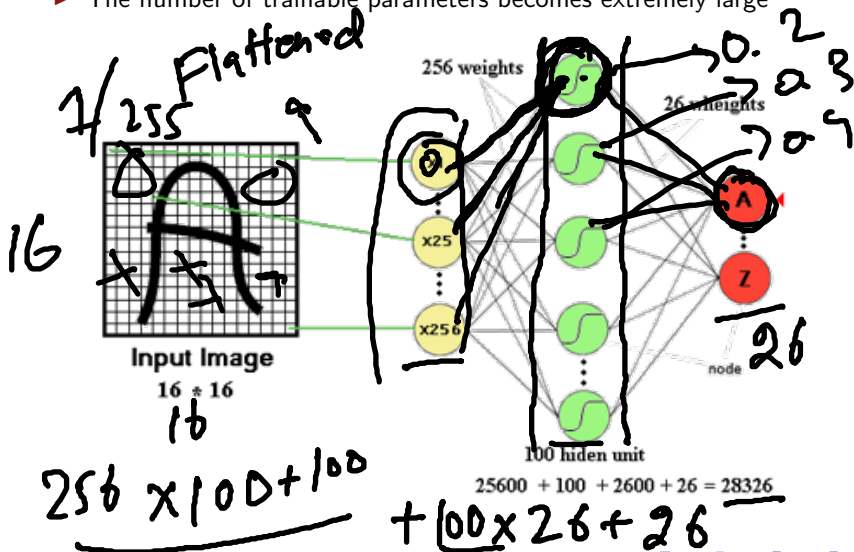
- ▶ Let's consider each pixel as a feature.



- ▶ This should work! But...

# Drawbacks of Fully-Connected Neural Networks

- ▶ The number of trainable parameters becomes extremely large



# Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ Little or no invariance to shifting, scaling, and other forms of distortion
- ▶ In other words, the Fully connected NN cannot recognize that the two images (original and shifted) represent the same object (letter "A").

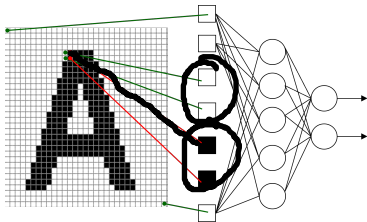


Figure 2: Original "A" character.

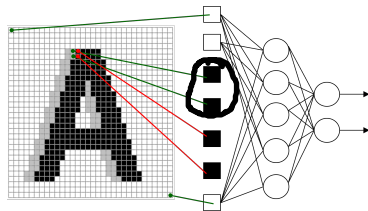


Figure 3: Shifted "A" character.

# Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ The topology of the input data is completely ignored (it treats pixels as independent features rather than structured patterns.)
- ▶ For a 32 × 32 image, we have
  - Black and white patterns:  $2^{\underline{32 \times 32}} = 2^{1024}$
  - Grayscale patterns 256<sup>32×32</sup> = 256<sup>1024</sup>

0, 1

2



- ▶ We need a model that can:
  - **Find Patterns in Images**: Recognize small features like edges and shapes in different parts of the image.
  - **Work Efficiently**: Avoid looking at every pixel separately by focusing on groups of pixels together.
  - **Handle Changes**: Still recognize the same object even if it moves, rotates, or looks slightly different.

- ▶ We need a model that can:
  - **Find Patterns in Images:** Recognize small features like edges and shapes in different parts of the image.
  - **Work Efficiently:** Avoid looking at every pixel separately by focusing on groups of pixels together.
  - **Handle Changes:** Still recognize the same object even if it moves, rotates, or looks slightly different.
  - We need **Convolutional Neural Networks (CNNs)**!

- ▶ CNNs are neural networks designed for image processing and consist of two main parts:
  - **Feature Extractor:** Automatically learns patterns (features) such as edges, textures, and shapes from the image.
  - **Classifier:** The extracted features are flattened into a vector and passed to a standard neural network (like the ones we used before) for classification.

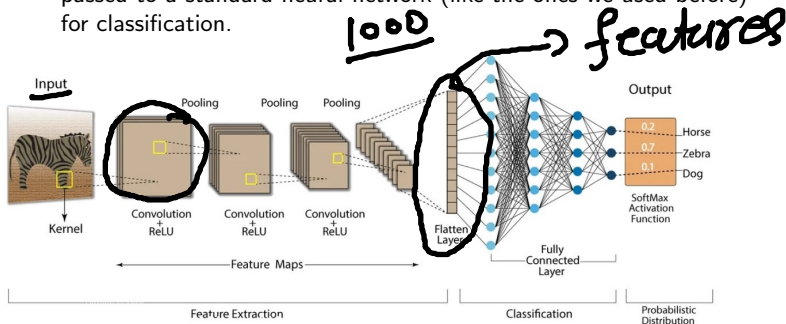


Figure 4: Convolutional Neural Network

- ▶ The feature extractor consists of three essential components:
  - **Convolution Layers:** Detect local patterns such as edges, textures, and shapes by applying filters to the image.
  - **Activation Function (ReLU):** Adds non-linearity.
  - **Pooling Layers:** Reduce the spatial size of extracted features.
- ▶ Let's start with the convolution layer.

- ▶ Let's consider this image.

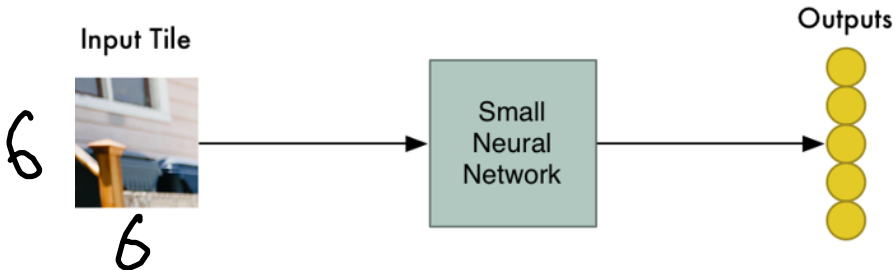


1. The image is divided into small overlapping tiles (regions).



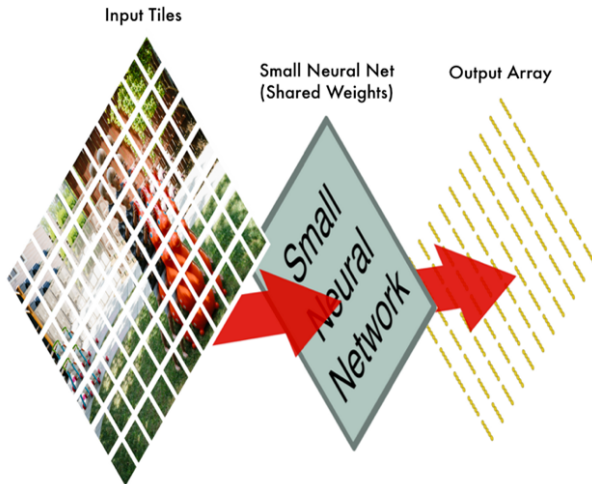
2. We process each tile using the weights matrix of a small neural network. We call this matrix a **kernel** or **filter**.

## Processing a single tile



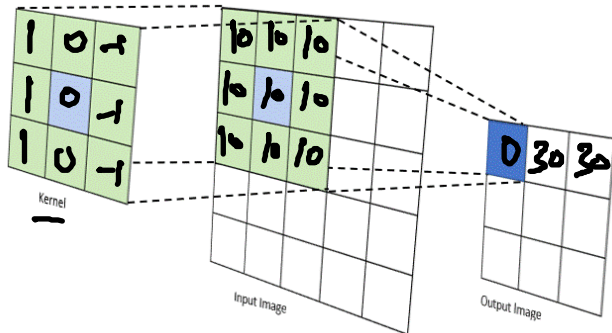
# How Convolution Works? (cont.)

3. Finally, the filter slides across the image (with the same weights) and processes all the tiles, creating an output **feature map**.



# How Convolution Works? (cont.)

- ▶ What we did in the last step is called the **Convolution Operation**.
- ▶ We convolved the kernel with the image, which means sliding the kernel over the entire image and computing the dot product between the kernel and small regions of the image at each step.



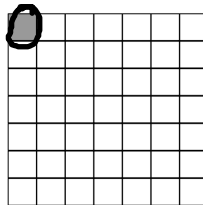
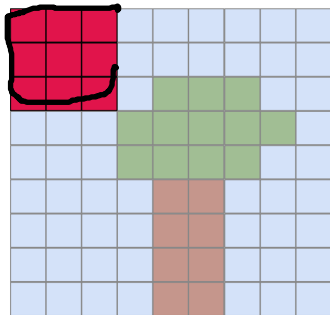
- ▶ This can be represented mathematically by:

$$z = W * x_{i,j} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} W_{ab} X_{(i+a)(j+b)}$$

- $z$ : Output of the convolution at  $(i, j)$ .
- $W$ : Filter (kernel) matrix of size  $m \times n$ .
- $x_{i,j}$ : Input value at position  $(i, j)$ .
- $W_{ab}$ : Weight of the filter at  $(a, b)$ .
- $X_{(i+a)(j+b)}$ : Input value in the small region at  $(i + a, j + b)$ .

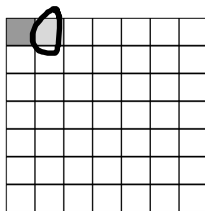
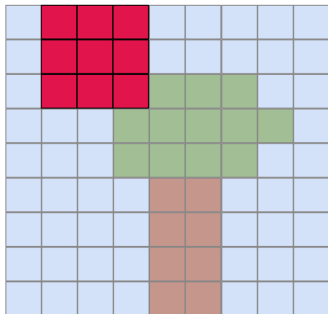
# How Convolution Works? (cont.)

3x3



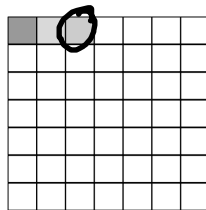
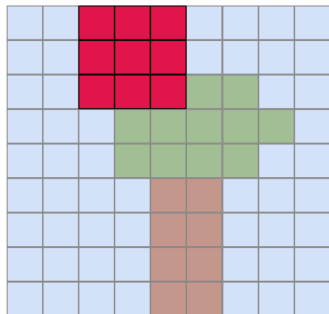
The **kernel** slides across the image and produces an output value at each position

# How Convolution Works? (cont.)



The **kernel** slides across the image and produces an output value at each position

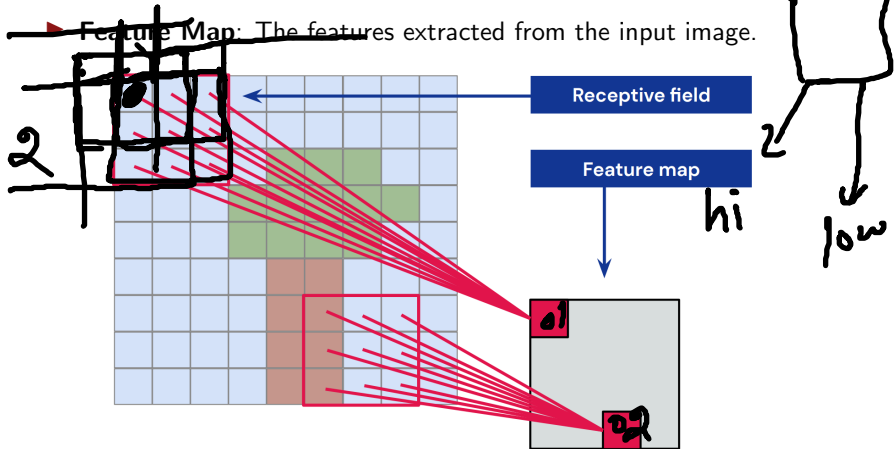
# How Convolution Works? (cont.)



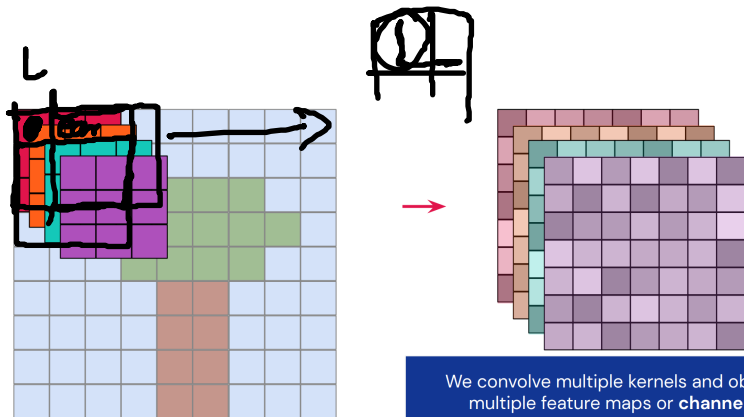
The **kernel** slides across the image and produces an output value at each position

# How Convolution Works? (cont.)

- ▶ **Receptive Field:** The region of the input image that the kernel operates on at each step.



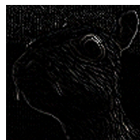
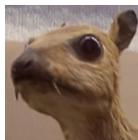
# How Convolution Works? (cont.)



- **Filters Example:** Different filters detect patterns like edges, textures, or smoothness, producing corresponding feature maps.

$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$$

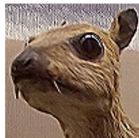
$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



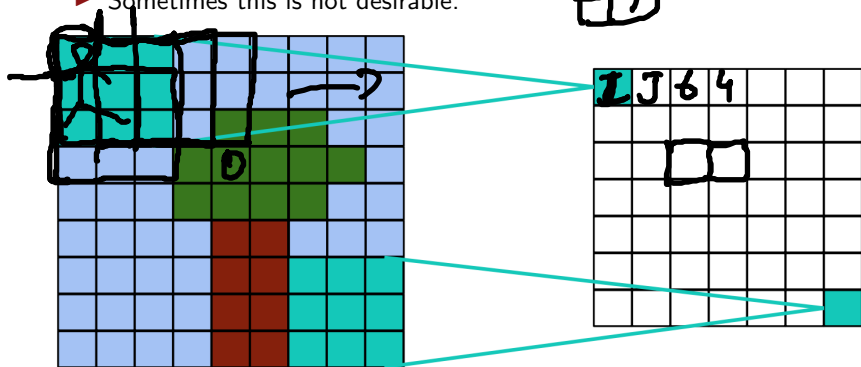
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



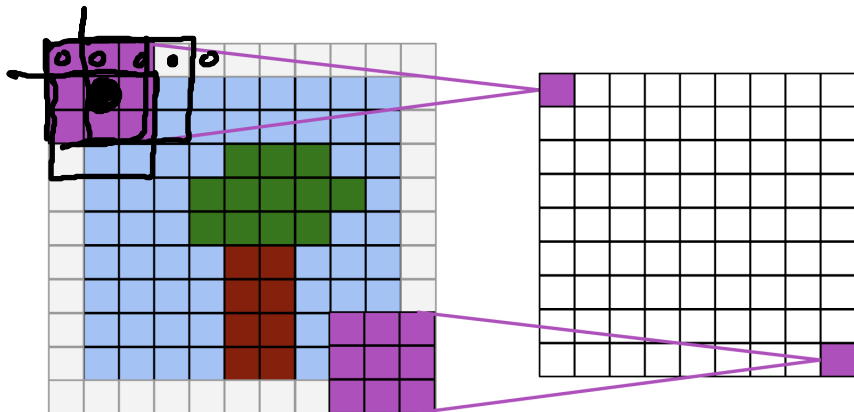
- ▶ Applying Convolution as such reduces the size of the borders.
- ▶ Sometimes this is not desirable.



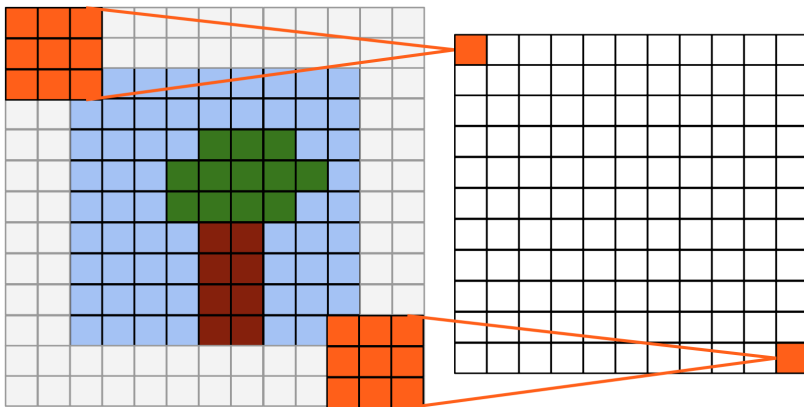
input size

# Controlling the Convolution Process (cont.)

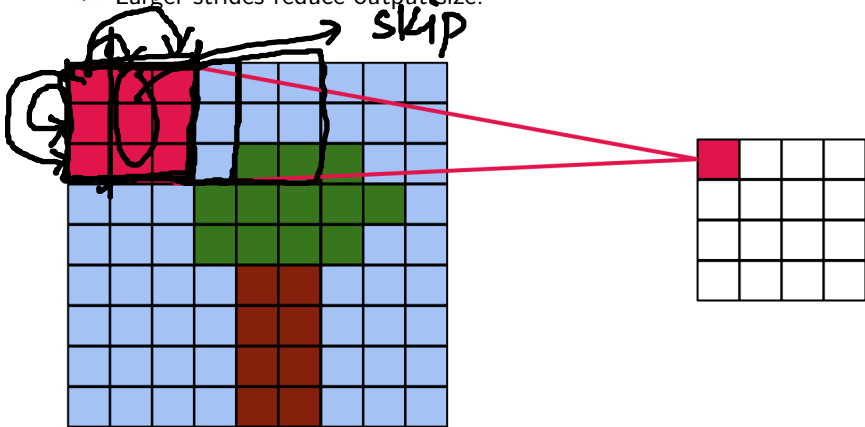
- ▶ Solution: We can use **Padding**, which pads the border with zeros.
- ▶ it has two types:
  1. **Same Convolution**: Padding is added so the output size equals the input size.



- ▶ 2. **Full Convolution:** Padding is added so the kernel covers the entire input, including edges.
- ▶ •  $\text{output size} = \text{input size} + \text{kernel size} - 1$

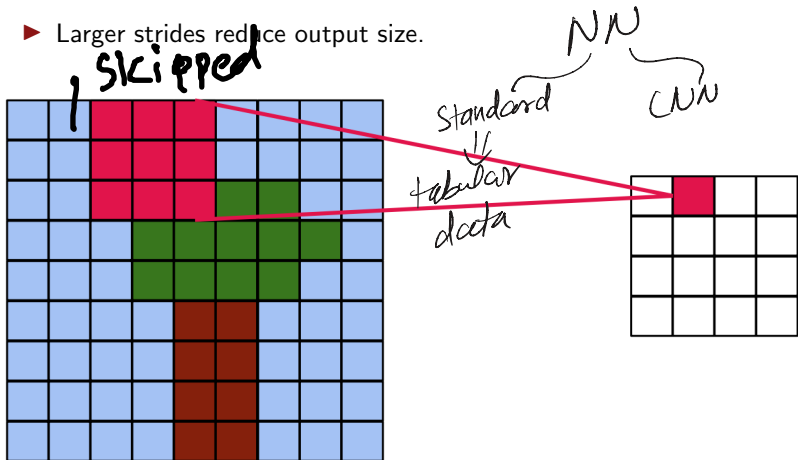


- ▶ **Strided Convolution:** Kernel slides along the image with a step  $> 1$
- ▶ Larger strides reduce output size.

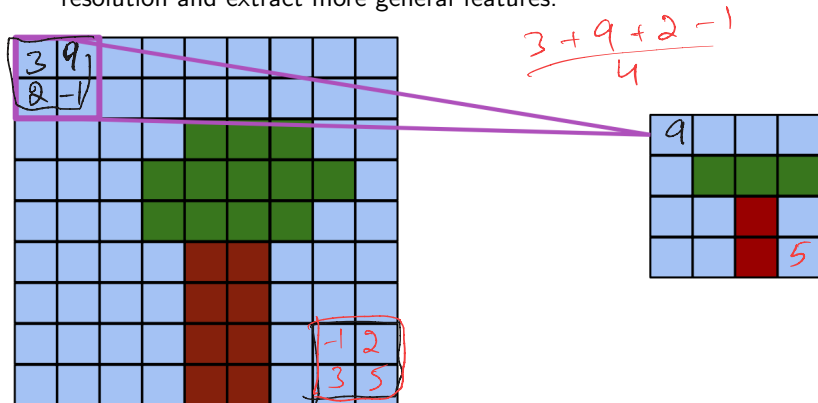


# Controlling the Convolution Process (cont.)

- ▶ **Strided Convolution:** Kernel slides along the image with a step  $> 1$
- ▶ Larger strides reduce output size.



- **Pooling:** Compute mean or max over small windows to reduce resolution and extract more general features.

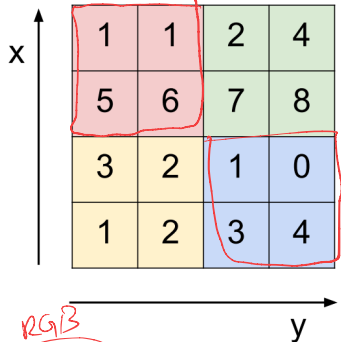


# Controlling the Convolution Process (cont.)

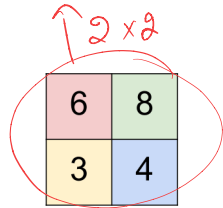
Group photo  
4x4

face

Single depth slice



max pool with 2x2 filters  
and stride 2

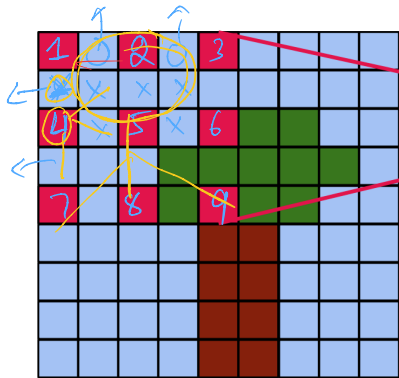


- No learnable parameters
- Introduces spatial invariance

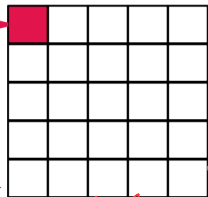


# Controlling the Convolution Process (cont.)

- ▶ **Dilated Convolution:** Kernel is spread out, step > 1 between kernel elements.
- ▶ It expands the receptive field without increasing the number of parameters, which makes it efficient.



3x3 filter/kernel



rate  
dilate

nn.Conv2D(

$S=2, P=1$

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

$n_{in}$  : number of input features

$n_{out}$  : number of output features

$k$  : convolution kernel size

$p$  : convolution padding size

$s$  : convolution stride size

- ▶ Just like Fully-Connected Neural Networks, we can apply an activation over convolutional layer outputs
- ▶ It helps break linearity
- ▶ For example, Rectified Linear Unit (ReLU):  $\sigma(x) = \max(0, x)$

Image  $\rightarrow$  Conv1  $\rightarrow$  features  
 $\downarrow$   
 Activat  
 $\downarrow$   
 Conv2

## Feature Map

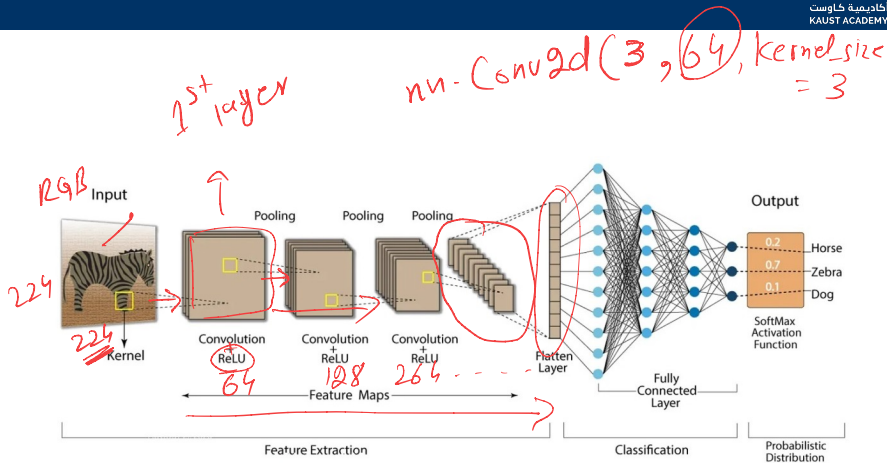
9	3	5	-8
-6	2	-3	1
1	3	4	1
3	-4	5	1

## ReLU Layer

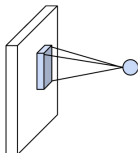


9	3	5	<del>0</del>
<del>0</del>	2	0	1
1	3	4	1
3	0	5	1

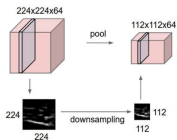
# Convolutional Neural Networks



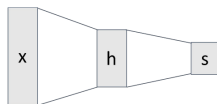
## Convolution Layers



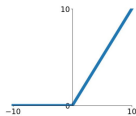
## Pooling Layers



## Fully-Connected Layers



## Activation Function



## Normalization

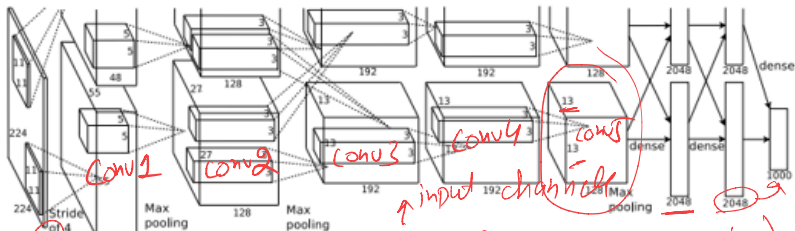
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

- ▶ Over time, researchers built advanced CNN architectures to improve performance and efficiency. These architectures introduced key innovations:
  - **AlexNet [Krizhevsky et al. 2012]**: The first CNN to achieve breakthrough performance on image classification.
  - **VGGNet [Simonyan and Zisserman, 2014]**: Used very deep networks (up to 19 layers).
  - **InceptionNet (GoogLeNet) [Szegedy et al., 2014]**: Used multiple filter sizes per layer (Inception modules).
  - **ResNet [He et al., 2015]**: Introduced skip connections for training very deep networks.
  - **EfficientNet [Tan and Le, 2019]**: Found a scaling method that simultaneously scales a CNN's depth, width, and resolution optimally using a single scaling coefficient.

- ▶ First big improvement in image classification.
- ▶ Made use of CNN, pooling, dropout, ReLU and training on GPUs.
- ▶ 5 convolutional layers, followed by max-pooling layers; with three fully connected layers at the end

$(\sigma, \max)$

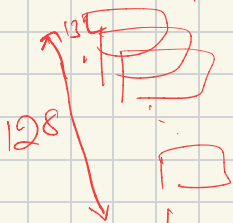
nn. ReLU(—)



Conv1 = nn.Conv2d(3, 48, k=11, s=4)

maxpool2 = nn.MaxPool2d(2) output channels

$$y_4 = \text{conv4}(y_3) \quad y_5 = \text{conv5}(y_4)$$

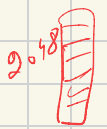


$$13 \times 13 \times 128 = 21632$$

conv1(x)

$$y_6 = \text{nn.Flatten}(y_5)$$

2048



$$\text{nn.Linear}(2048, 2048)$$

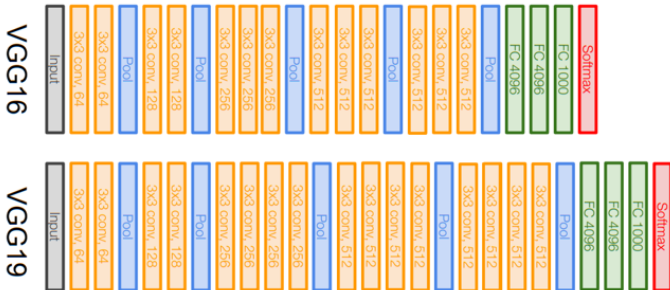


$$\text{nn.Linear}(2048, 1000)$$

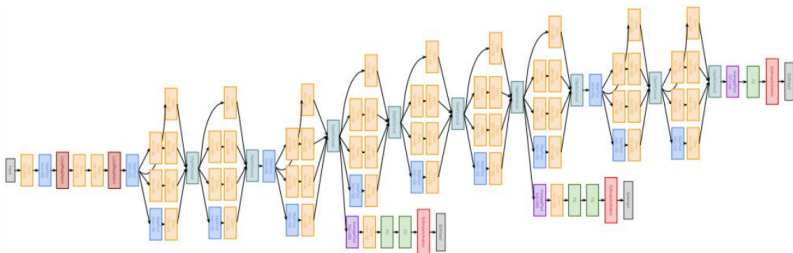


cat → highest value  
dog → low  
bus → low

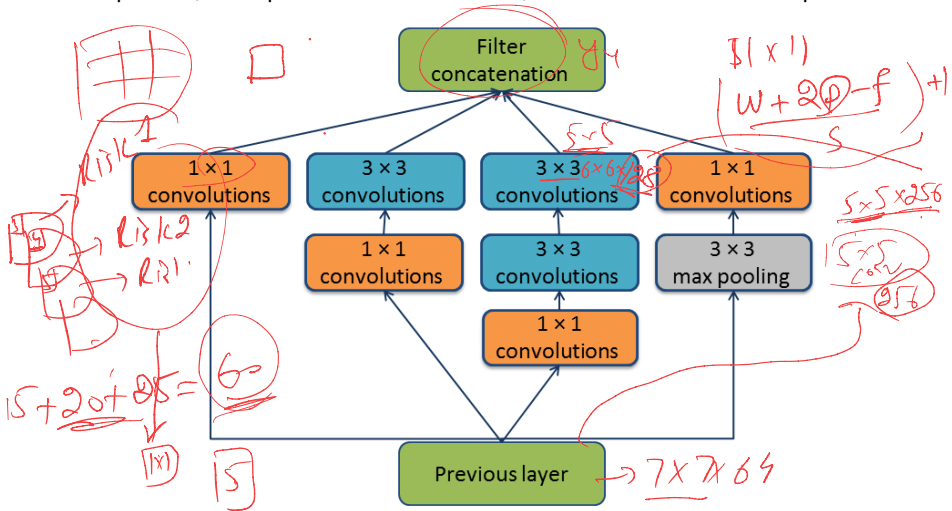
- ▶ **Improvement over AlexNet:** Uses a deeper network with small filters instead of a shallow network with larger filters.
- ▶ A stack of 3x3 conv layers (vs. 11x11 conv) has same receptive field, more non-linearities, fewer parameters and deeper network representation.
- ▶ Two variants: VGG16 or VGG19 conv layers plus 3 FC layers.

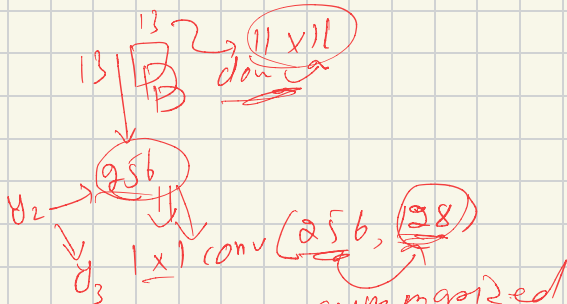
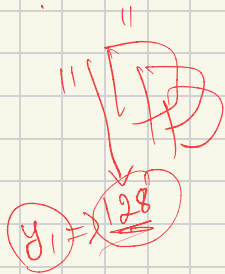


- ▶ Going Deep: 22 layers
- ▶ Only 5 million parameters! ( $12\times$  fewer than AlexNet,  $27\times$  fewer than VGGNet).
- ▶ Introduced efficient "Inception module"
- ▶ Introduced "bottleneck" layers that use  $1\times 1$  convolutions to reduce the number of channels and mix information across them.



- **Inception module:** Uses multiple filter sizes ( $1 \times 1$ ,  $3 \times 3$ ,  $5 \times 5$ ), in parallel, to capture different features, then combines their outputs.

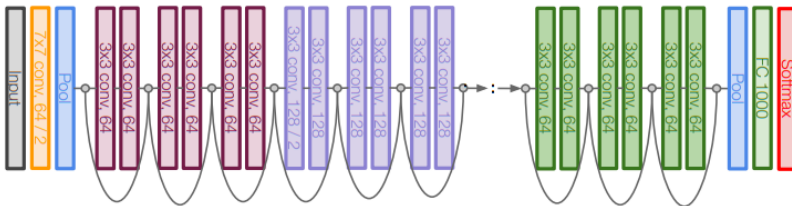




$$y_4 = \text{torch.cat}(y_1, y_3) \text{ summarized}$$

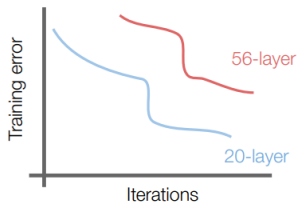
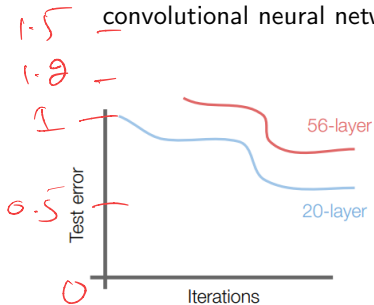
- ▶ **Problem:** Making networks deeper does not always improve accuracy.
- ▶ **Why?** In very deep networks, gradients become extremely small as they move backward through layers, making learning slow or stopping it altogether (**vanishing gradient problem**).
- ▶ **Solution:** Residual Network (ResNet) introduces **skip connections (residuals)**, allowing information to flow more easily.

- ▶ Very deep networks using residual connections
- ▶ 152-layer model for ImageNet
- ▶ Stacked Residual Blocks
- ▶ **Residual:** A shortcut connection that helps the network pass information through layers more easily.

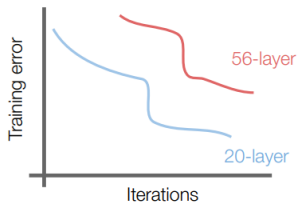
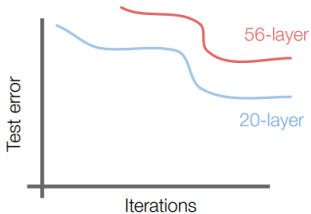


- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

- What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

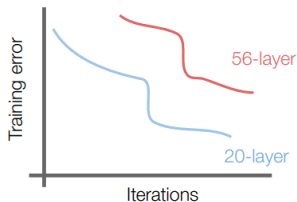
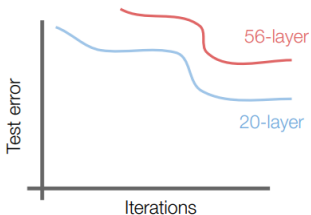


- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



- ▶ 56-layer model performs worse on both test and training error

- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



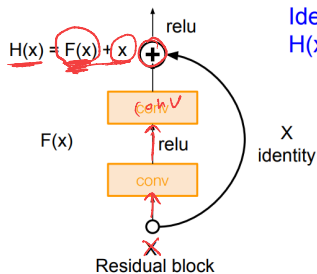
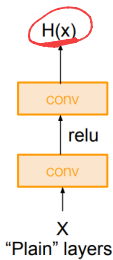
- ▶ 56-layer model performs worse on both test and training error
- ▶ The deeper model performs worse, but it's not caused by overfitting!

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize
- ▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize
- ▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

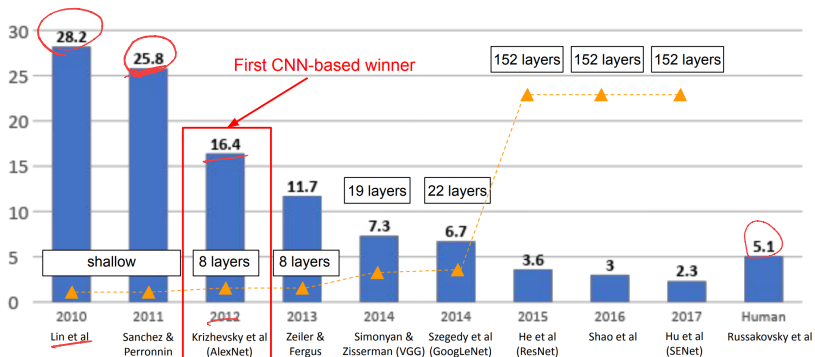


Identity mapping:  
 $H(x) = x$  if  $F(x) = 0$

- ▶ The most extensive data for Image Classification
- ▶ 3 RGB channels from 0 to 255
- ▶ 14,197,122 images
- ▶ 1000 classes

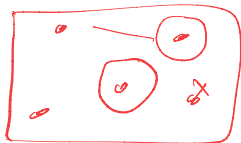


# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



► **Problem:** To improve accuracy, we can:

- Increase the number of layers (**depth**)
- Increase the number of neurons in each layer (**width**)
- Use higher-resolution images (**resolution**)



But finding the right balance of these three was largely based on trial and error.

► **Solution:** EfficientNet introduced a **compound scaling** method—a mathematical formula to systematically find the optimal balance among **depth**, **width**, and **resolution**.

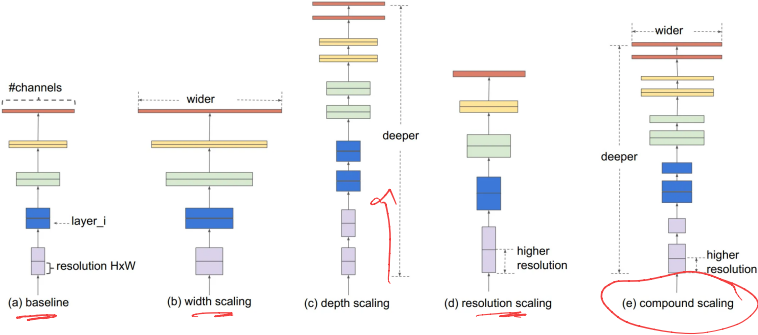


Figure 5: EfficientNet Scaling.

- ▶ **Compound Scaling** Uses a single scaling coefficient ( $\phi$ ) to control:
  - **Network Depth** ( $\alpha^\phi$ ) → More layers
  - **Network Width** ( $\beta^\phi$ ) → More channels per layer
  - **Input Resolution** ( $\gamma^\phi$ ) → Larger input images
- ▶ The goal: find  $\alpha, \beta, \gamma$  that balance accuracy & efficiency, then scale up optimally by increasing the global coefficient  $\phi$ .

- ▶ EfficientNet optimizes depth, width, and resolution using this constraint:

$$\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

*floating point operations.*

- ▶ **Why this equation?**

20% • Increasing **depth** ( $\alpha$ ) increases FLOPs **linearly**.

10% • Increasing **width** ( $\beta$ ) increases FLOPs **quadratically** ( $\beta^2$ ).

15% • Increasing **resolution** ( $\gamma$ ) increases FLOPs **quadratically** ( $\gamma^2$ ).

- ▶ To double total FLOPs, the three factors must be balanced together.

0.3 x 0.5

- ▶ The authors of EfficientNet searched for the best scaling factors on a small baseline model.
- ▶ They found:

$$\begin{array}{l} \alpha = \underline{1.2}, \quad \beta = \underline{1.1}, \quad \gamma = \underline{1.15} \\ \alpha = 1 + \frac{0.2}{\downarrow} \left| \beta = 1 + 0.1 \right. \left. \delta = 1 + 0.15 \right. \\ 20\% \qquad \qquad \qquad 10\% \qquad \qquad \qquad 15\% \end{array}$$

- ▶ The EfficientNet family (B0 to B7) is generated using:

$$\text{Depth} = 1.2^\phi, \quad \text{Width} = 1.1^\phi, \quad \text{Resolution} = 1.15^\phi$$

Model	$\phi$	Depth ( $\alpha^\phi$ )	Width ( $\beta^\phi$ )
B0	0	$1.2^0 = 1.0$	$1.1^0 = 1.0$
B1	1	$1.2^1 = 1.2$	$1.1^1 = 1.1$
B2	2	$1.2^2 = 1.44$	$1.1^2 = 1.21$
B3	3	$1.2^3 = 1.73$	$1.1^3 = 1.33$
B4	4	$1.2^4 = 2.07$	$1.1^4 = 1.46$
B5	5	$1.2^5 = 2.49$	$1.1^5 = 1.61$
B6	6	$1.2^6 = 2.99$	$1.1^6 = 1.77$
B7	7	$1.2^7 = 3.58$	$1.1^7 = 1.94$

Table 1: Scaling EfficientNet from B0 to B7

- ▶ We multiply these scaling factors by the baseline EfficientNet-B0 values and round to the nearest integer to get the new depth, width, and resolution for each model.

- ▶ EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.

- ▶ EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.
- ▶ EfficientNet-B7 reaches 84.4% Top-1 and 97.3% Top-5 accuracy on ImageNet.

- ▶ EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.
- ▶ EfficientNet-B7 reaches 84.4% Top-1 and 97.3% Top-5 accuracy on ImageNet.
- ▶ More efficient than previous CNN models—8.4x smaller and 6.1x faster than competitors.

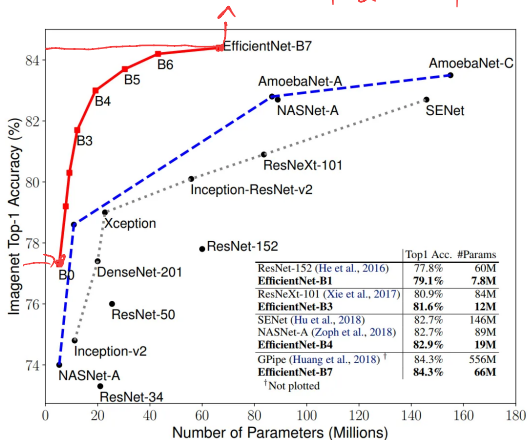


Figure 6: EfficientNet Performance on Imagenet.

## ► Why MobileNets?

- Small-sized models are crucial for mobile and embedded devices.
- MobileNets reduce computational cost and memory usage while maintaining good accuracy.

## ► Key Idea:

- Use **depthwise-separable convolutions** to significantly reduce computation compared to standard convolutions.

- ▶ **Computational cost of standard convolution:**

$$\text{Cost} = \# \text{ filter params} \times \# \text{ filter positions} \times \# \text{ filters}$$

- ▶ Filters operate on all input channels, increasing computation significantly.



$$X = \left[ \begin{array}{cccccc} 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & & & & & \\ \vdots & & & & & \\ \vdots & & & & & \\ \vdots & & & & & \end{array} \right] \begin{array}{l} \rightarrow R \\ \rightarrow G \\ \rightarrow B \end{array}$$

$$3 \times 3 \times \underbrace{(3)}_{RGB}$$

$$6 \times 6 \times 3$$

$$(K_n \times K_w \times \underline{C_{in+1}}) \underline{C_{out}}$$

$$(3 \times 3 \times 3 + 1) \times 10 = (27 + 1) \times 10 = \underline{280}$$

$$\left[ \begin{array}{c} \left( \begin{array}{cc} 1 & 0 \\ 0 & -1 \end{array} \right)^R \\ \left( \begin{array}{c} 1 \\ 0 \\ -1 \end{array} \right)^G \\ \left( \begin{array}{c} 1 \\ 0 \\ -1 \end{array} \right)^B \end{array} \right]$$

$$\left( \begin{matrix} \left( \begin{matrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{matrix} \right)^R \\ \left( \begin{matrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{matrix} \right)^G \\ \left( \begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} \right)^B \end{matrix} \right) * \left( \begin{matrix} \left( \begin{matrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{matrix} \right)^R \\ \left( \begin{matrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{matrix} \right)^G \\ \left( \begin{matrix} -1 & 0 \\ 0 & 0 \\ 0 & 0 \end{matrix} \right)^B \end{matrix} \right)$$

$$\begin{aligned}
 R & 30 + 0 + 0 = 30 \\
 G & = 30 + \\
 B & = 30 = \mathbf{90}
 \end{aligned}$$

Depth-wise separable convolution

$$\textcircled{1} \left( \begin{matrix} \left( \begin{matrix} 10 & 10 & 0 \\ 10 & 10 & 0 \\ 10 & 10 & 0 \end{matrix} \right)^R \\ \left( \begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} \right)^G \\ \left( \begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} \right)^B \end{matrix} \right) * \left( \begin{matrix} \left( \begin{matrix} 1 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \right)^R \\ \left( \begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} \right)^G \\ \left( \begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} \right)^B \end{matrix} \right)$$

$$\left( \begin{matrix} R \\ 30 \\ G \\ 30 \\ B \\ 30 \end{matrix} \right)$$

$$\left( \begin{matrix} 2 \\ \end{matrix} \right)^B$$

Point wise convolution

$$\begin{matrix} R & & G & & B \\ \left[ \begin{array}{ccc} 1 & & 1 \\ & 1 & \\ & & 1 \end{array} \right] & \times & \left[ \begin{array}{ccc} 30 & 30 & 30 \end{array} \right] \\ & & 30 & + & 30 & + & 30 & = & 90 \end{matrix}$$

$$\textcircled{a} \quad \left[ \underline{3 \times 3} \times \textcircled{3} = 27 + 3 = \textcircled{30} \right] = \boxed{70}$$

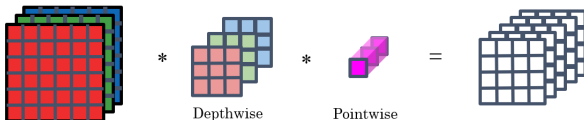
$$\textcircled{b} \quad \left[ \underline{1 \times 1} \times \textcircled{3} \right] \underline{10} = \underline{30} + 10 = \textcircled{40}$$

- ▶ Split standard convolution into two steps:
  - **Depthwise Convolution:** Applies a single filter per input channel.
  - **Pointwise Convolution:** Combines outputs from depthwise convolution.
- ▶ **Key Benefit:** Reduces computational cost significantly compared to standard convolution.

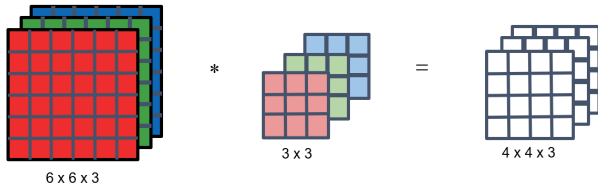
Normal Convolution



Depthwise Separable Convolution

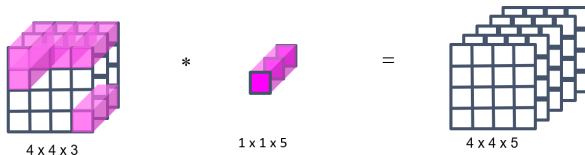


- ▶ Operates on each input channel separately.



Computational cost = #filter params x # filter positions x # of filters

- Combines outputs from depthwise convolution using  $1 \times 1$  convolutions (mixes channels).



$$\text{Computational cost} = \# \text{filter params} \times \# \text{filter positions} \times \# \text{of filters}$$

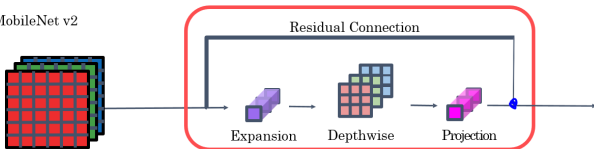
## ► MobileNet v2:

- Adds **residual connections**.
- Introduces:
  - **Expansion step:** Expands input dimensions before depthwise convolution.
  - **Projection step:** Reduces dimensions after processing.

MobileNet v1



MobileNet v2

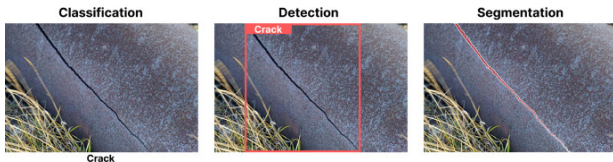


## ► CNN Architectures Used:

- ResNet
- VGG
- InceptionNet

## ► Use Case:

- Automated analysis of pipeline conditions using camera feeds.
- Classification of corrosion severity levels.
- Predictive maintenance through visual pattern recognition.



## ► Benefits:

- Reduces manual inspection costs and time.
- Early detection prevents catastrophic failures.
- Enables data-driven maintenance scheduling.

## ► Links & Resources:

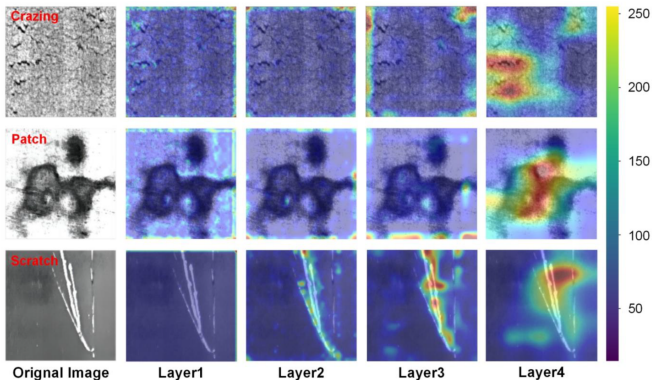
- [viso.ai](https://viso.ai)
- [ScienceDirect](https://www.sciencedirect.com)
- [NextBrain](https://www.nextbrain.com)

## ▶ CNN Architectures Used:

- ResNet50 with Transfer Learning
- VGG19
- AlexNet

## ▶ Use Case:

- Classify types of defects in steel surfaces and equipment.
- Automated quality control for manufactured components.
- Multi-class defect recognition (rust, cracks, deformation, etc.).



## ► Benefits:

- 99.4% accuracy with transfer learning approaches.
- Cost savings through early detection.
- Replaces subjective manual inspections.

## ▶ Links & Resources:

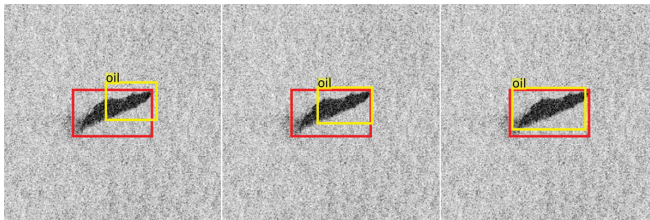
- MDPI: Steel defect detection (ResNet, 99.4% accuracy)

## ▶ CNN Architectures Used:

- YOLOv4
- Mask R-CNN
- CNNs

## ▶ Use Case:

- Classify oil spill severity and type from SAR/optical imagery.
- Automated environmental monitoring systems.
- Multi-class classification for response prioritization.



(a) IoU = 30%

(b) IoU = 45%

(c) IoU = 75%

## ► Benefits:

- Faster emergency response decisions.
- Accurate spill characterization for cleanup planning.
- Compliance with environmental regulations.

## ► Research Papers & Resources:

- [Deep learning pink spill detector \(Sentinel-1 SAR, 2022\)](#)
- [Comprehensive review of oil spill DL methods \(2024\)](#)

These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: [Deep Learning for Computer Vision](#)
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV [Deep Learning for Computer Vision: Fundamentals and Applications](#)
- ▶ Justin Johnson, UMich EECS 498.008/598.008: [Deep Learning for Computer Vision](#)
- ▶ Sander Dieleman, Deepmind: [Deep Learning Lecture Series 2020](#)